

(19) 日本国特許庁 (J P)

(12) 公開特許公報 (A)

(11) 特許出願公開番号

特開平11-306205

(43) 公開日 平成11年(1999)11月5日

(51) IntCl.⁶

G 0 6 F 17/30

識別記号

F I

G 0 6 F 15/403
15/40

3 3 0 C
3 7 0 A

審査請求 有 請求項の数 7 F D (全 21 頁)

(21) 出願番号 特願平10-129485

(22) 出願日 平成10年(1998)4月23日

(71) 出願人 000004237

日本電気株式会社

東京都港区芝五丁目7番1号

(72) 発明者 島津 秀雄

東京都港区芝五丁目7番1号 日本電気株式会社内

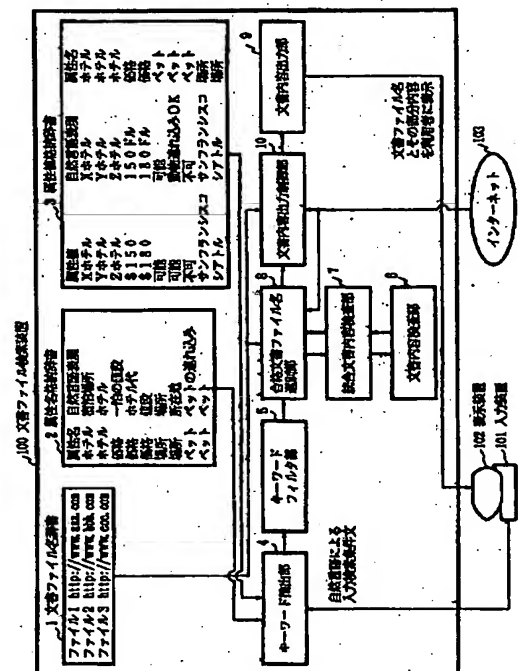
(74) 代理人 弁理士 境 廣巳

(54) 【発明の名称】 文書ファイル検索装置及びプログラムを記録した機械読み取り可能な記録媒体

(57) 【要約】

【課題】 WWW のホームページに対する自然言語による検索問い合わせを実現する。

【解決手段】 検索対象文書ファイルであるWWW のホームページをXML で記述する。検索条件文が入力されるとキーワード抽出部4は、属性名を表現する自然言語表現はその属性名を含む属性名インデックスに、属性値を表現する自然言語表現はその属性値と属性名との対を含む属性値インデックスに変換する。キーワードフィルタ部5は、変換後のインデックス列中で、同一の属性の属性名インデックスと属性値インデックスとが隣どうしに存在する箇所の属性名インデックスを削除する。文書内容検査部6は、検索対象文書ファイル中に、変換後のインデックス列中の全ての属性値インデックスの属性名と属性値との対に対応するタグが存在するか否かを調べ、存在する場合、文書内容出力部9が変換後のインデックス列中の属性名インデックスの属性名を持つタグの属性値を検索して出力する。



1

【特許請求の範囲】

【請求項 1】 属性の属性名とその属性の属性値との対を内蔵する文書ファイルを検索対象文書ファイルとし、検索対象文書ファイルから、利用者が自然言語で指定した検索条件に適合する部分を検索する文書ファイル検索装置において、

自然言語で表現した検索要求文を先頭から順に探査し、属性名を表現する自然言語表現に対してはその属性名を属性名インデックスとして出力し、属性値を表現する自然言語表現に対してはその属性値と属性名との対を属性値インデックスとして出力することを順次行うキーワード抽出部と、

前記キーワード抽出部の出力を入力して先頭から順に探査し、同一の属性の属性名インデックスと属性値インデックスとが隣どうしに存在する場合のみ、前記属性名インデックスを削除し、それ以外の部分はそのまま出力するキーワードフィルタ部と、

検索対象文書ファイル中に、前記キーワードフィルタ部から出力された全ての属性値インデックスの属性名と属性値との対が内蔵されているか否かを調べ、内蔵されている場合、前記キーワードフィルタ部から出力された属性名インデックスの属性名に対応する属性値を検索対象文書ファイルから検索して出力する検索手段とを備えた文書ファイル検索装置。

【請求項 2】 検索対象文書ファイル中に存在する属性名について、属性名とその属性名を表現する自然言語表現との対を格納しておく属性名格納辞書と、

検索対象文書ファイル中に存在する属性値について、属性値とその属性値に対応する属性名とその属性値を表現する自然言語表現との 3 つ組を格納しておく属性値格納辞書とを備え、

前記キーワード抽出部は、自然言語で表現した検索要求文を先頭から順に探査し、属性名格納辞書を参照して、属性名を表現する自然言語表現が含まれていたら、その自然言語表現と対である属性名を属性名インデックスとして出力し、属性値格納辞書を参照して、属性値を表現する自然言語表現が含まれていたら、その自然言語表現と 3 つ組である属性値と属性名との対の集合を属性値インデックスとして出力する構成を有することを特徴とする請求項 1 記載の文書ファイル検索装置。

【請求項 3】 文書中に書かれた意味を表現する属性名のついたタグとその属性の属性値との対を内蔵する文書ファイルを検索対象文書ファイルとし、検索対象文書ファイルから、利用者が自然言語で指定した検索条件に適合する部分を検索する文書ファイル検索装置において、検索対象文書ファイル中に存在する属性名について、属性名とその属性名を表現する自然言語表現との対の集合を格納しておく属性名格納辞書と、

検索対象文書ファイル中に存在する属性値について、属性値とその属性値に対応する属性名とその属性値を表現

2

する自然言語表現との 3 つ組の集合を格納しておく属性値格納辞書と、

自然言語で表現した検索要求文を先頭から順に探査し、属性名格納辞書を参照して、属性名を表現する自然言語表現が含まれていたら、その自然言語表現と対である属性名を属性名インデックスとして出力し、属性値格納辞書を参照して、属性値を表現する自然言語表現が含まれていたら、その自然言語表現と 3 つ組である属性値と属性名との対を属性値インデックスとして出力するキーワード抽出部と、

キーワード抽出部の出力を入力し、先頭から順に探査し、同一の属性の属性名インデックスと属性値インデックスとが隣どうしに存在する場合のみ、前記属性名インデックスを削除し、それ以外の部分はそのまま出力するキーワードフィルタ部と、

検索対象文書ファイル中に、前記キーワードフィルタ部から出力された全ての属性値インデックスの属性名と属性値との対に対応するタグの対が内蔵されているか否かを調べ、内蔵されている場合、前記キーワードフィルタ部から出力された属性名インデックスの属性名を持つタグの属性値を検索対象文書ファイルから検索して出力する検索手段とを備えた文書ファイル検索装置。

【請求項 4】 文書中に書かれた意味を表現する属性名のついたタグとその属性の値との対を複数個内蔵する文書ファイルの集合から、利用者が自然言語で指定した検索条件を満足する文書ファイルを選択してその適合する部分を表示する文書ファイル検索装置において、検索対象となるすべての文書ファイルの名前と存在位置とを格納する文書ファイル名辞書と、

検索対象となる文書ファイル中に存在する属性名について、属性名とその属性名を表現する自然言語表現との対の集合を格納しておく属性名格納辞書と、

検索対象となる文書ファイル中に存在する属性値について、属性値とその属性値に対応する属性名とその属性値を表現する自然言語表現との 3 つ組の集合を格納しておく属性値格納辞書と、

利用者が、自然言語で表現した検索要求文を入力すると、前記入力文を先頭から順に探査し、属性名格納辞書を参照して、属性名を表現する自然言語表現が含まれていたら、その自然言語表現と対である属性名を属性名インデックスとして出力し、属性値格納辞書を参照して、属性値を表現する自然言語表現が含まれていたら、その自然言語表現と 3 つ組である属性値と属性名との対の集合を属性値インデックスとして出力することを順次行うキーワード抽出部と、

キーワード抽出部の出力を入力し、先頭から順に探査し、同一の属性の属性名インデックスと属性値インデックスとが隣通しに存在する場合のみ、前記属性名インデックスを削除し、それ以外の部分はそのまま出力するキーワードフィルタ部と、

10

20

30

40

50

3

文書ファイルの内容と属性値インデックスとを入力すると、前記文書ファイルの内容中に、前記属性値インデックス中の属性名を含むタグが存在するかどうか調べ、存在する場合は、そのタグと対で存在する属性値を取り出し、その値が前記属性値インデックス中の属性値と等しいかどうか調べ、等しい場合は、合格の出力をし、そうでない場合は不合格の出力をする文書内容検査部と、文書ファイルの内容と1つ以上の属性値インデックスとを入力すると、前記属性値インデックスから1つずつ取り出し、前記文書ファイルの内容と前記取り出した属性値インデックスとを1つずつ文書内容検査部に渡していき、すべての属性値インデックスに対してその出力が合格のときは、合格を出力し、そうでないときは不合格を出力する統合文書内容検査部と、前記文書ファイル名辞書を参照して、1つずつ文書ファイルの内容を取り出し、前記文書の内容とキーワードフィルタ部の出力のうちの属性値インデックスの部分とを統合文書内容検査部に渡し、前記統合文書内容検査部の出力を受け取ることを前記1つずつ取り出した文書ファイルのすべてに対して行い、前記出力が合格の文書ファイルの名前のみを出力する合格文書ファイル名選別部と、

文書ファイル名と前記文書ファイル名の内容とキーワードフィルタ部の出力である属性名インデックスとを入力すると、前記属性名インデックスのうちの1つを取り出し、与えられた前記文書ファイルの内容中に、前記取り出した属性名を含むタグが存在するかどうか調べ、存在する場合は、その属性名のタグの値と前記入力した文書ファイル名とを利用者に表示し、存在しない場合には何も出力しないことを、前記入力した属性名インデックスのそれぞれに対して行う文書内容出力部と、前記合格文書ファイル名選別部の出力である文書ファイル名の集合を入力し、文書ファイル名格納辞書を参照して、前記入力した文書ファイル名の集合の要素を1つずつ取り出し、文書内容出力部に渡すことを、前記入力中の文書ファイル名のすべてに対して行うことを繰り返す文書内容出力制御部とを備えることを特徴とする文書ファイル検索装置。

【請求項5】 属性の属性名とその属性の属性値との対を内蔵する文書ファイルを検索対象文書ファイルとし、検索対象文書ファイルから、利用者が自然言語で指定した検索条件に適合する部分を検索する文書ファイル検索装置を構成するコンピュータを、

自然言語で表現した検索要求文を先頭から順に探査し、属性名を表現する自然言語表現に対してはその属性名を属性名インデックスとして出力し、属性値を表現する自然言語表現に対してはその属性値と属性名との対を属性値インデックスとして出力することを順次行うキーワード抽出部、

前記キーワード抽出部の出力を入力して先頭から順に探

4

査し、同一の属性の属性名インデックスと属性値インデックスとが隣どうしに存在する場合のみ、前記属性名インデックスを削除し、それ以外の部分はそのまま出力するキーワードフィルタ部、

検索対象文書ファイル中に、前記キーワードフィルタ部から出力された全ての属性値インデックスの属性名と属性値との対が内蔵されているか否かを調べ、内蔵されている場合、前記キーワードフィルタ部から出力された属性名インデックスの属性名に対応する属性値を検索対象文書ファイルから検索して出力する検索手段、

として機能させるプログラムを記録した機械読み取り可能な記録媒体。

【請求項6】 文書中に書かれた意味を表現する属性名のついたタグとその属性の属性値との対を内蔵する文書ファイルを検索対象文書ファイルとし、検索対象文書ファイルから、利用者が自然言語で指定した検索条件に適合する部分を検索する文書ファイル検索装置を構成するコンピュータを、

検索対象文書ファイル中に存在する属性名について、属性名とその属性名を表現する自然言語表現との対の集合を格納しておく属性名格納辞書、

検索対象文書ファイル中に存在する属性値について、属性値とその属性値に対応する属性名とその属性値を表現する自然言語表現との3つ組の集合を格納しておく属性値格納辞書、

自然言語で表現した検索要求文を先頭から順に探査し、属性名格納辞書を参照して、属性名を表現する自然言語表現が含まれていたら、その自然言語表現と対である属性名を属性名インデックスとして出力し、属性値格納辞書を参照して、属性値を表現する自然言語表現が含まれていたら、その自然言語表現と3つ組である属性値と属性名との対を属性値インデックスとして出力するキーワード抽出部、

キーワード抽出部の出力を入力し、先頭から順に探査し、同一の属性の属性名インデックスと属性値インデックスとが隣どうしに存在する場合のみ、前記属性名インデックスを削除し、それ以外の部分はそのまま出力するキーワードフィルタ部、

検索対象文書ファイル中に、前記キーワードフィルタ部から出力された全ての属性値インデックスの属性名と属性値との対に対応するタグの対が内蔵されているか否かを調べ、内蔵されている場合、前記キーワードフィルタ部から出力された属性名インデックスの属性名を持つタグの属性値を検索対象文書ファイルから検索して出力する検索手段、

として機能させるプログラムを記録した機械読み取り可能な記録媒体。

【請求項7】 文書中に書かれた意味を表現する属性名のついたタグとその属性の値との対を複数個内蔵する文書ファイルの集合から、利用者が自然言語で指定した検

10

20

30

40

50

5

索条件を満足する文書ファイルを選択してその適合する部分を表示する文書ファイル検索装置を構成するコンピュータを、

検索対象となるすべての文書ファイルの名前と存在位置とを格納する文書ファイル名辞書、

検索対象となる文書ファイル中に存在する属性名について、属性名とその属性名を表現する自然言語表現との対の集合を格納しておく属性名格納辞書、

検索対象となる文書ファイル中に存在する属性値について、属性値とその属性値に対応する属性名とその属性値を表現する自然言語表現との3つ組の集合を格納しておく属性値格納辞書、

利用者が、自然言語で表現した検索要求文を入力すると、前記入力文を先頭から順に探査し、属性名格納辞書を参照して、属性名を表現する自然言語表現が含まれていたなら、その自然言語表現と対である属性名を属性名インデックスとして出力し、属性値格納辞書を参照して、属性値を表現する自然言語表現が含まれていたなら、その自然言語表現と3つ組である属性値と属性名との対の集合を属性値インデックスとして出力することを順次行うキーワード抽出部、

キーワード抽出部の出力を入力し、先頭から順に探査し、同一の属性の属性名インデックスと属性値インデックスとが隣通しに存在する場合のみ、前記属性名インデックスを削除し、それ以外の部分はそのまま出力するキーワードフィルタ部、

文書ファイルの内容と属性値インデックスとを入力すると、前記文書ファイルの内容中に、前記属性値インデックス中の属性名を含むタグが存在するかどうか調べ、存在する場合は、そのタグと対で存在する属性値を取り出し、その値が前記属性値インデックス中の属性値と等しいかどうか調べ、等しい場合は、合格の出力をし、そうでない場合は不合格の出力をする文書内容検査部、

文書ファイルの内容と1つ以上の属性値インデックスを入力すると、前記属性値インデックスから1つずつ取り出し、前記文書ファイルの内容と前記取り出した属性値インデックスとを1つずつ文書内容検査部に渡していき、すべての属性値インデックスに対してその出力が合格のときは、合格を出力し、そうでないときは不合格を出力する統合文書内容検査部、

前記文書ファイル名辞書を参照して、1つずつ文書ファイルの内容を取り出し、前記文書の内容とキーワードフィルタ部の出力のうちの属性値インデックスの部分とを統合文書内容検査部に渡し、前記統合文書内容検査部の出力を受け取ることを前記1つずつ取り出した文書ファイルのすべてに対して行い、前記出力が合格の文書ファイルの名前のみを出力する合格文書ファイル名選別部、文書ファイル名と前記文書ファイル名の内容とキーワードフィルタ部の出力である属性名インデックスとを入力すると、前記属性名インデックスのうちの1つを取り出

6

し、与えられた前記文書ファイルの内容中に、前記取り出した属性名を含むタグが存在するかどうか調べ、存在する場合は、その属性名のタグの値と前記入力した文書ファイル名とを利用者に表示し、存在しない場合には何も出力しないことを、前記入力した属性名インデックスのそれぞれに対して行う文書内容出力部、

前記合格文書ファイル名選別部の出力である文書ファイル名の集合を入力し、文書ファイル名格納辞書を参照して、前記入力した文書ファイル名の集合の要素を1つずつ取り出し、文書内容出力部に渡すことを、前記入力中の文書ファイル名のすべてに対して行うことを繰り返す文書内容出力制御部、

として機能させるプログラムを記録した機械読み取り可能な記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は文書ファイル検索装置に関し、特に自然言語による検索問い合わせを可能とした文書ファイル検索装置に関する。

【0002】

【従来の技術】一般に情報検索において利用者の検索意図をより精密に表現させようとする場合には、日本語や英語のような自然言語によってそれを表現させる方法が有効である。データベースに対する検索を自然言語で行うシステムは既に存在し、自然言語インタフェースと呼ばれている(参考文献:ディベロップングアナチュラランゲージインタフェースツコンプレックスデータ、ジー・ジー・ヘンドリックス他, "Developing a Natural Language Interface to Complex Data", ACM Trans. on Database Systems, 1978.)。

【0003】従来の自然言語インタフェースは、利用者の自然言語による検索問い合わせを解釈して、その問い合わせをデータベースの検索言語(SQL)の検索式に変換し、その検索式をデータベースシステムに送り、データベースシステムから戻された検索結果を利用者に提示するものである。

【0004】しかし、従来の自然言語インタフェースは、既に商用化が始まって20年以上たったのにも関わらず、まだ実用のレベルに達していない。その理由の1つは、自然言語インタフェースシステムが利用者の自由な問い合わせを解釈することができず、そのシステムが許容する構文や語彙が明確に限定されているので、結局は利用者はどういう言い回しを使えるかを覚えなくてはならないためである。従って、自然言語インタフェースといっても複雑なコマンド体系と変わらない。つまり、従来の自然言語インタフェースは、利用者の自由な言い回しを受け付けることが出来ないと言うことが問題であった(参考文献:ディベロップングアナチュラランゲージインタフェースツコンプレックスデータ、ジー・ジー・ヘンドリックス他, "Developing a Natural Language

age Interface to Complex Data", ACM Trans. on Database Systems, 1978.)。とりわけ、簡単な質問については、利用者は完璧な自然言語でなく、簡潔な言い回しをを使いたいの、それが許されないということが問題であった。簡潔な言い回しの例としては、キーワード列による表現や非文法的な言い回し、あるいは自然言語文の一部、等がある。

【0005】そこで本出願人は、先の特許出願において、簡潔な言い回しの自然言語による検索を実現する「自然言語解釈方法」を提案した（特開平5-67136号公報）。これは、属性の属性名とその属性の属性値との組の集まりであるデータベースのテーブルを検索対象とし、自然言語による問い合わせ文中の各単語を属性名と属性値とその他とに分類し、属性名に分類された単語は応答属性名として保存し、属性値に分類された単語はその属性値とそれに対応する属性名とを組にして条件属性値組群として保存し、この保存された条件属性値組群中の属性値と属性名との組が全てテーブル中に存在する場合、前記応答属性名として保存された各属性名に対

150 ドル、サンフランシスコ、ホテル、ベット可能
とキーワードを並べたら、非常に大量のホームページのリストが出力される。

【0007】自然言語インタフェースシステムを導入しWWWのホームページを検索できれば、上の例のような精密な検索条件を素直に表現し、適切なホームページのみが検索できることになる。しかしながら、WWW上のホームページの検索に対し自然言語インタフェースを適用した例は見当たらない。

【0008】なお、WWWのホームページに対する検索技術の他の例として、特開平10-40262号公報に記載された「情報検索装置」があるが、これは、感性表現データをキーワードにした検索を可能にすることで、明確な検索対象または検索条件を持たない利用者の感性に合った情報検索を目的としているため、上の例のような精密な検索条件による検索には向いていない。

【0009】

【発明が解決しようとする課題】上述した特開平5-67136号公報に記載された技術は、単純な方法で自然言語による問い合わせを解釈することができるので、自然言語インタフェースシステムを実用化する上で有効な手段となり得るが、未だ解決すべき課題が残されている。それは、自然言語による問い合わせ文中に或る属性名が存在する場合、それが必ず応答属性名として扱われ、問い合わせに対する回答中に含められるため、回答が冗長になる場合があることである。

【0010】例えば、「属性名＝書名、その属性値＝人間失格、属性名＝著者、その属性値＝太宰治」を持つテーブルに対して、「書名＝人間失格の著者は？」という問い合わせを行った場合、まず、「書名」が属性名と判定されて応答属性名として保存され、次いで「人間失

*応する前記テーブル中の属性値を、問い合わせに対する回答として出力するものである。なお、これに類似する従来技術として、やはり本出願人によって先に出願された特開平5-242147号公報にかかる「自然言語解釈方法」がある。

【0006】他方、最近におけるWorld Wide Web (WWW)の利用の拡大に伴い、WWW上での検索技術の重要性が高まっている。WWWの利用者がWWW上で情報検索をするときに使う典型的なツールはサーチエンジンである。この例としては、Altavista、Infoseek、Lycosなどがある。しかし、サーチエンジンでは、キーワードを組み合わせる検索の形式なので、利用者の検索の意図が直接的に反映させられないことが多い。例えば、ホテルに関する情報を探すときに、値段が150ドルでベットを連れて行くのが可能で、しかもサンフランシスコ近辺にあるホテルのホームページを見つけないときに、そのような検索意図をキーワードの並びのみで表現することは不可能である。仮に、(1)式で表現するように

格」が属性値と判定されて属性値「人間失格」とその属性名である「書名」との組が条件属性値組群として保存され、次いで、「著者」が属性名と判定されて応答属性名として保存される。そして、条件属性値組群中の属性値「人間失格」と属性名「書名」との組を有する前記テーブルが検索され、そのテーブルから応答属性名「書名」と「著者」とに対応する属性値「人間失格」と「太宰治」とが検索されて出力される。つまり、「人間失格」をも出力している分、回答が冗長になっている。

【0011】そこで本発明の目的は、自然言語による検索問い合わせに対する回答の冗長性を極力無くすことにある。

【0012】また、本発明の他の目的は、WWWのホームページに対しても自然言語による検索問い合わせを可能にすることにある。

【0013】

【課題を解決するための手段】(1)第1の発明
上述した特開平5-67136号公報に記載された技術において、回答が冗長になっている理由は、自然言語による問い合わせ文中に属性名が存在する場合、それに対応する属性値を利用者が問い合わせ文中で記述しているにもかかわらず、一律に応答属性名として扱っているためである。そこで、本発明では、自然言語による問い合わせ文中で検索条件を指定するために或る属性名とそれに対応する属性値とを記述する場合、利用者はそれらを互いに隣接して記述する傾向にある点に着目し、同一の属性の属性名と属性値とが隣どうしに現れる場合にその属性名を応答属性名に含めないようにしている。より具体的には、属性の属性名とその属性の属性値との対を内蔵する文書ファイルを検索対象文書ファイルとし、検索

対象文書ファイルから、利用者が自然言語で指定した検索条件に適合する部分を検索する文書ファイル検索装置において、自然言語で表現した検索要求文を先頭から順に探査し、属性名を表現する自然言語表現に対してはその属性名を属性名インデックスとして出力し、属性値を表現する自然言語表現に対してはその属性値と属性名との対を属性値インデックスとして出力することを順次行うキーワード抽出部と、前記キーワード抽出部の出力を入力して先頭から順に探査し、同一の属性の属性名インデックスと属性値インデックスとが隣どうしに存在する場合のみ、前記属性名インデックスを削除し、それ以外の部分はそのまま出力するキーワードフィルタ部と、検索対象文書ファイル中に、前記キーワードフィルタ部から出力された全ての属性値インデックスの属性名と属性値との対が内蔵されているか否かを調べ、内蔵されている場合、前記キーワードフィルタ部から出力された属性名インデックスの属性名に対応する属性値を検索対象文書ファイルから検索して出力する検索手段とを備えている。

【0014】更に、検索要求文中のどの自然言語表現が属性名を表現し、またどの自然言語表現が属性値を表現しているかを正確に判定できるようにするために、検索対象文書ファイル中に存在する属性名について、属性名とその属性名を表現する自然言語表現との対を格納しておく属性名格納辞書と、検索対象文書ファイル中に存在する属性値について、属性値とその属性値に対応する属性名とその属性値を表現する自然言語表現との3つ組を格納しておく属性値格納辞書とを備え、前記キーワード抽出部は、自然言語で表現した検索要求文を先頭から順に探査し、属性名格納辞書を参照して、属性名を表現する自然言語表現が含まれていたなら、その自然言語表現と対である属性名を属性名インデックスとして出力し、属性値格納辞書を参照して、属性値を表現する自然言語表現が含まれていたなら、その自然言語表現と3つ組である属性値と属性名との対の集合を属性値インデックスとして出力する構成を有する。

「値段=\$150 かつ ベット=可能 かつ 場所=サンフランシスコ」

… (2)

のように変換できた。これは、SQL 言語にそのまま変換される。しかし、WWW のホームページは、通常このような属性名と属性値の情報が入っていないので、SQL の式で表現できるような検索式に変換できない。つまり、従来のWWW のホームページ作成言語はHTMLである（参考文献：ワールドワイドウェブコンソーシアムのホームページ、URL <http://www.w3.org>）。HTMLでは、文書ファイル中に、その構成を表現するための属性名と属性値の組が内蔵されている。例えば、図6に示すのが、HTMLファイルの例である。ここで、<と>に囲まれたものが属性タグであり、単独で使われるもの（例：）と、開始タグ（例：<TR>）と終了タグ（例：</TR>）の対で

*【0015】このように構成された本発明の文書ファイル検索装置にあっては、利用者が自然言語で表現した検索要求文を入力すると、まずキーワード抽出部が、検索要求文を先頭から順に探査し、属性名を表現する自然言語表現に対してはその属性名を属性名インデックスとして出力し、属性値を表現する自然言語表現に対してはその属性値と属性名との対を属性値インデックスとして出力し、次いでキーワードフィルタ部が、キーワード抽出部の出力を入力して先頭から順に探査し、同一の属性の属性名インデックスと属性値インデックスとが隣どうしに存在する場合のみ、属性名インデックスを削除し、次いで、検索手段が、検索対象文書ファイル中に、キーワードフィルタ部から出力された全ての属性値インデックスの属性名と属性値との対が内蔵されているか否かを調べ、内蔵されている場合、キーワードフィルタ部から出力された属性名インデックスの属性名に対応する属性値を検索対象文書ファイルから検索して出力することにより、利用者への回答が冗長になるのを防いでいる。

【0016】(2) 第2の発明

WWW 上のホームページの検索に対し自然言語インタフェースを適用するのが困難であった理由は、WWW のホームページの中身が自然言語で書かれた文章や図から構成されるファイルであり、データベースのように、属性名と属性値の集合でないことである。つまり、従来の自然言語インタフェースシステムが対象とするデータベースは、属性名と属性値の集合だったため、従来の技術の項の例で出てくるホテルのデータベースがあるとする、（名前：Xホテル、値段：\$150、ベット：可能、場所：サンフランシスコ）
（名前：Yホテル、値段：\$200、ベット：不可、場所：ロサンゼルス）
（名前：Zホテル、値段：\$180、ベット：不可、場所：シアトル）
のような形態で格納されており、利用者の問い合わせは、(2) 式で表現するように

使われるものがある。HTMLのタグの特徴は、それが文書ファイル中の外見の表現を定義するのに限定されていることである。例えば、表的な表現にするタグは<TABLE>であり、改行を表すタグは<P>で表現される。このようなHTMLファイルをWWW ブラウザに読み込ませると、図7に示すような形態になってユーザに出力表示される。しかし、HTMLでは、文書中の意味を表現する為のタグを定義することは出来ない。

【0017】そこで本発明では、WWW の文書ファイルに、その文書ファイル中の意味を表現する属性名と属性値との組を内蔵させる。具体的には、例えば、ファイル中に文書の内容を属性タグとその属性値の対の集合で表

現できるように拡張したXML (Extensible Markup Language)で文書を記述する(参考文献:ワールドワイドウェブコンソーシアムのホームページ、「エクステンシブルマークアップ ランゲージ 1.0」<http://www.w3.org/TR/PR-xml-971208>)。XML は、WWW の標準を決める機関であるワールドワイドウェブコンソーシアム(参考文献:ワールドワイドウェブコンソーシアムのホームページ、URL <http://www.w3.org>)によって1997年12月にその仕様が提案された。XML で記述された文書では、文書の内容を機械が可読になって内容による検索が可能になる。そこで、本発明ではそのことを利用してWWW のホームページに対して自然言語による検索問い合わせを実現する。

【0018】具体的には、文書中に書かれた意味を表現する属性名のついたタグとその属性の属性値との対を内蔵する文書ファイルを検索対象文書ファイルとし、検索対象文書ファイルから、利用者が自然言語で指定した検索条件に適合する部分を検索する文書ファイル検索装置において、検索対象文書ファイル中に存在する属性名について、属性名とその属性名を表現する自然言語表現との対の集合を格納しておく属性名格納辞書と、検索対象文書ファイル中に存在する属性値について、属性値とその属性値に対応する属性名とその属性値を表現する自然言語表現との3つ組の集合を格納しておく属性値格納辞書と、自然言語で表現した検索要求文を先頭から順に探査し、属性名格納辞書を参照して、属性名を表現する自然言語表現が含まれていたなら、その自然言語表現と対である属性名を属性名インデックスとして出力し、属性値格納辞書を参照して、属性値を表現する自然言語表現が含まれていたなら、その自然言語表現と3つ組である属性値と属性名との対の集合を属性値インデックスとして出力するキーワード抽出部と、キーワード抽出部の出力を入力し、先頭から順に探査し、同一の属性の属性名インデックスと属性値インデックスとが隣どうしに存在する場合のみ、前記属性名インデックスを削除し、それ以外の部分はそのまま出力するキーワードフィルタ部と、検索対象文書ファイル中に、前記キーワードフィルタ部から出力された全ての属性値インデックスの属性名と属性値との対に対応するタグが内蔵されているか否かを調べ、内蔵されている場合、前記キーワードフィルタ部から出力された属性名インデックスの属性名を持つタグの属性値を検索対象文書ファイルから検索して出力する検索手段とを備えている。

【0019】このように構成された本発明の文書ファイル検索装置にあっては、利用者が自然言語で指定した検索要求文を入力すると、キーワード抽出部が、検索要求文を先頭から順に探査し、属性名を表現する自然言語表現が含まれていたなら、その属性名を属性名インデックスとして出力し、属性値を表現する自然言語表現が含まれていたなら、その属性値と属性名との対の集合を属性値イ

ンデックスとして出力し、次いで、キーワードフィルタ部が、キーワード抽出部の出力を入力し、先頭から順に探査し、同一の属性の属性名インデックスと属性値インデックスとが隣どうしに存在する場合のみ、前記属性名インデックスを削除し、次いで、検索手段が、検索対象文書ファイル中に、キーワードフィルタ部から出力された全ての属性値インデックスの属性名と属性値との対に対応するタグが内蔵されているか否かを調べ、内蔵されている場合、キーワードフィルタ部から出力された属性名インデックスの属性名を持つタグの属性値を検索対象文書ファイルから検索して出力する。

【0020】また、予め登録された多数の文書ファイルの中から利用者が自然言語で入力した検索条件を満たす文書ファイルのみを選別し、さらにその中の利用者が必要な部分を利用者に表示できるようにするために、文書中に書かれた意味を表現する属性名のついたタグとその属性の値との対を複数個内蔵する文書ファイルの集合から、利用者が自然言語で指定した検索条件を満足する文書ファイルを選択してその適合する部分を表示する文書ファイル検索装置において、検索対象となるすべての文書ファイルの名前と存在位置とを格納する文書ファイル名辞書と、検索対象となる文書ファイル中に存在する属性名について、属性名とその属性名を表現する自然言語表現との対の集合を格納しておく属性名格納辞書と、検索対象となる文書ファイル中に存在する属性値について、属性値とその属性値に対応する属性名とその属性値を表現する自然言語表現との3つ組の集合を格納しておく属性値格納辞書と、利用者が、自然言語で表現した検索要求文を入力すると、前記入力文を先頭から順に探査し、属性名格納辞書を参照して、属性名を表現する自然言語表現が含まれていたなら、その自然言語表現と対である属性名を属性名インデックスとして出力し、属性値格納辞書を参照して、属性値を表現する自然言語表現が含まれていたなら、その自然言語表現と3つ組である属性値と属性名との対の集合を属性値インデックスとして出力することを順次行うキーワード抽出部と、キーワード抽出部の出力を入力し、先頭から順に探査し、同一の属性の属性名インデックスと属性値インデックスとが隣どうしに存在する場合のみ、前記属性名インデックスを削除し、それ以外の部分はそのまま出力するキーワードフィルタ部と、文書ファイルの内容と属性値インデックスとを入力すると、前記文書ファイルの内容中に、前記属性値インデックス中の属性名を含むタグが存在するかどうか調べ、存在する場合は、そのタグと対で存在する属性値を取り出し、その値が前記属性値インデックス中の属性値と等しいかどうか調べ、等しい場合は、合格の出力をし、そうでない場合は不合格の出力をする文書内容検査部と、文書ファイルの内容と1つ以上の属性値インデックスとを入力すると、前記属性値インデックスから1つずつ取り出し、前記文書ファイルの内容と前記取り出

13

した属性値インデックスを1つずつ文書内容検査部に渡していき、すべての属性値インデックスに対してその出力が合格のときは、合格を出力し、そうでないときは不合格を出力する統合文書内容検査部と、文書ファイル名辞書を参照して、1つずつ文書ファイルの内容を取り出し、前記文書の内容とキーワードフィルタ部の出力のうちの属性値インデックスの部分とを統合文書内容検査部に渡し、前記統合文書内容検査部の出力を受け取ること
 を前記1つずつ取り出した文書ファイルのすべてに対して行い、前記出力が合格の文書ファイルの名前のみを出力する合格文書ファイル名選別部と、文書ファイル名と前記文書ファイル名の内容とキーワードフィルタ部の出力である属性名インデックスとを入力すると、前記属性名インデックスのうちの1つを取り出し、与えられた前記文書ファイルの内容中に、前記取り出した属性名を含むタグが存在するかどうか調べ、存在する場合は、その属性名のタグの値と前記入力した文書ファイル名とを利用者に表示し、存在しない場合には何も出力しないことを、前記入力した属性名インデックスのそれぞれに対して行う文書内容出力部と、前記合格文書ファイル名選別部の出力である文書ファイル名の集合を入力し、文書ファイル名格納辞書を参照して、前記入力した文書ファイル名格納辞書を参照して、前記入力した文書ファイル名から生成されたキーワード列：

価格（「価格」の属性名インデックス）、
 \$150ドル（「価格」の属性値インデックス）、
 ベット（「ベット可能性」の属性名インデックス）、
 可能（「ベット可能性」の属性値インデックス）、
 サンフランシスコ（「場所」の属性値インデックス）、
 ホテル（「ホテル名」の属性名インデックス）

…(4)

【0024】次に、属性名インデックスと属性値インデックスの並び順を参照して、冗長な部分の統合を行う。
 同一の属性に対する属性名インデックスと属性値インデックスとが隣りどうしに並んでいるときには、属性名インデックスの方を削除する。上のキーワード列例は、次のように圧縮される。

圧縮されたキーワード列：

\$150ドル（「価格」の属性値インデックス）、
 可能（「ベット可能性」の属性値インデックス）、
 サンフランシスコ（「場所」の属性値インデックス）、
 ホテル（「ホテル名」の属性名インデックス）

…(5)

【0025】次に、抽出したキーワード列を解釈する。
 属性値インデックスは、それが参照する属性の値として、属性値インデックスが保持する値を取ること、という条件式と解釈する。例えば、
 \$150ドル（「価格」の属性値インデックス）
 は、

全体の条件式

「「価格」属性の値 = \$150」 かつ 50

14

*ル名の集合の要素を1つずつ取り出し、文書内容出力部に渡すことを、前記入力中の文書ファイル名のすべてに対して行うことを繰り返す文書内容出力制御部とを備えている。

【0021】このように構成された本発明の文書ファイル検索装置の作用を、その理解を容易にするために、例を使って説明する。まず、利用者が検索する対象となるWWW文書ファイルとして、図5(a)、(b)に示したものを使用する。図5の文書ファイル中には、文章テキストの他に、属性の属性名とその属性の属性値との対が含まれている。また、利用者の検索文の例として、次の文を使う。

検索入力文：「値段が150ドルでベットを連れて行くのが可能で、しかも、サンフランシスコ近辺にあるホテルの情報を見つけない」

【0022】まず、第1段階では、入力文をキーワード列に変換する。キーワードの種類としては、2種類存在する。1つ目は、属性名を参照する自然言語表現であり、属性名インデックスと呼ぶ。2つ目は、属性値を参照する自然言語表現であり、属性値インデックスと呼ぶ。

【0023】

※「「価格」属性の値 = \$150」という解釈をする。

★「「価格」属性の値 = \$150」という解釈をする。

【0026】複数の属性値インデックスが存在する場合は、それらの解釈を論理積したものが全体の条件式となる。上の例では、以下になる。

15

「「ペット可能性」属性の値 = 可能」 かつ
 「「場所」属性の値 = サンフランシスコ」

16

…(6)

【0027】属性名インデックスは、それが参照する属性の値を出力せよ、という解釈になる。上の例では、以 *
 検索部分の特定

ホテル（「ホテル名」の属性名インデックス）」

…(7)

【0028】この意味は、「「ホテル名」属性の値を出力せよ」という解釈となる。複数の属性名インデックスがあるときは、それら複数の属性名インデックスを順次出力せよ、という意味になる。

【0029】入力文全体の解釈は、属性値インデックスから生成される検索条件式を満足するWWW上の文書ファイルを選択し、次に、それらの文書ファイル中から属性名インデックスの解釈で指定される属性名の値を抽出してそれを利用者に表示すれば良い。

【0030】

【発明の実施の形態】図1を参照すると、本発明の実施の形態の文書ファイル検索装置100は、文書ファイル名辞書1と、属性名格納辞書2と、属性値格納辞書3と、キーワード抽出部4と、キーワードフィルタ部5と、文書内容検査部6と、統合文書内容検査部7と、合格文書ファイル名選別部8と、文書内容出力部9と、文書内容出力制御部10とから構成され、キーボード等の入力装置101、CRTディスプレイ等の表示装置102およびインターネット103に接続されている。

【0031】文書ファイル名辞書1には、検索対象となるすべての文書ファイルの名前とその物理的な位置とが格納されている。検索対象となる文書ファイルがHTMLやXMLで記述されている場合には、文書ファイルは、世界中のWWWサーバに分散していることも可能である。その場合、文書ファイルの位置は、「http://.....」というURL記述になる。

【0032】属性名格納辞書2には、検索対象となる文書ファイル中に存在する属性タグの属性名とその属性名を表現する自然言語表現との対が登録されている。ある属性名を参照する自然言語表現の中の最も基本的なものは、その属性名そのものである。例えば、「ホテル」という属性名を参照する自然言語表現としては、「ホテ

ル」である。しかし、それ以外にも、「ホテル」を参照する表現がある。例えば、「宿泊場所」、「泊まる場所」などの表現がある。これらが、下記の表1で示すような対になって登録される。

【0033】

【表1】

属性名	自然言語表現
ホテル	ホテル
ホテル	宿泊場所
ホテル	泊まる場所

【0034】属性値格納辞書3には、検索対象となる文書ファイル中に存在する属性値について、属性値とその属性値に対応する属性名とその属性値を表現する自然言語表現との3つ組が格納される。ある属性値を参照する自然言語表現としてもっとも基本的なものは、その属性値そのものである。例えば、「Xホテル」という属性値を参照する自然言語表現としては、「Xホテル」そのものがありこれ以外にはないかもしれない。しかし、別の例では、「ペット」属性の属性値を表わす自然言語表現としては、「可能」の他に「動物連れ込みOK」「ペット同伴OK」「犬猫可」のような表現も登録しておいてもよい。属性値格納辞書3には、下記の表2で示すように3つ組でデータが格納される。

【0035】

【表2】

属性値	自然言語表現	対応する属性名
可能	可能	ベット
可能	動物連れ込みOK	ベット
可能	ペット同伴OK	ベット
可能	犬猫可	ベット
不可	動物不可	ベット

【0036】キーワード抽出部4は、自然言語表現による入力条件検索文を入力装置101を通じて利用者から受け取ると、属性名格納辞書2と属性値格納辞書3とを参照して、その中の自然言語表現として登録されている表現が入力条件検索文中にないかどうかを調べる。あった場合には、それが属性名の場合には、属性名のみを出力する。この出力のことを属性名インデックスと呼ぶ。他方、それが属性値の場合には、属性値と対応する属性名との対を出力する。この出力のことを属性値インデックスと呼ぶ。これらは、入力条件検索文の先頭から調べていき、マッチするものが見つかったら、その順番に出力していく。

【0037】キーワードフィルタ部5は、キーワード抽出部4の出力をそのまま受け取り、先頭から順に探索し、同一の属性の属性名インデックスと属性値インデックスとが隣どうしに存在する場合は、その属性名インデックスを削除し、それ以外の部分はそのまま素通しで出力する。

【0038】文書内容検査部6は、統合文書内容検査部7から文書ファイルの内容である文字列と属性値インデックスとを入力として受け付ける。入力として受け取った文書ファイルの内容文字列中に、受け取った属性値インデックス中の属性名を含むタグが存在するかどうか調べ、存在する場合は、そのタグと対で存在する属性値を取り出し、その値がこの属性値インデックス中の属性値と等しいかどうか調べ等しい場合は、合格の出力をし、そうでない場合は不合格の出力をする。文書内容検査部6は、統合文書内容検査部7から呼び出されて動作する一種のサブルーチンの役割を果たしている。

【0039】統合文書内容検査部7は、合格文書ファイル名選別部8から文書ファイルの内容である文字列と1つ以上の属性値インデックスとを入力として受け付ける。与えられた属性値インデックスは、1つ1つが「属性値インデックス中に記述された属性の値として、属性値インデックス中に記述された値をとらねばならない」という条件式を表現していると見做す。統合文書内容検査部7の役割は、与えられた文字列中から、属性値イン

デックスに記述された属性表現を見つけて、その条件が満足されているかを調べることである。入力として与えられた1つ以上の属性値インデックスのすべての条件を満足すれば、「合格」という値を出力し、そうでない場合は、「不合格」という値を出力する。実際に、文書ファイルの内容である文字列が1つの属性値インデックスの条件を満足するかどうかを判定するのは、文書内容検査部6が行う。統合文書内容検査部7は、複数の属性値インデックスがあった場合に、属性値インデックス1つずつを文書内容検査部6に順々に渡していく一種のループ制御を行っている。統合文書内容検査部7も、合格文書ファイル名選別部8から呼び出されるサブルーチンの役割である。

【0040】合格文書ファイル名選別部8は、文書ファイル名辞書1を参照して、必要に応じてインターネット103を通じて世界中に分散しているWWWサーバをアクセスして1つずつ文書ファイルの内容を取り出し、この文書の内容とキーワードフィルタ部5の出力のうち属性値インデックスの部分とを統合文書内容検査部7に渡し、統合文書内容検査部7の出力を受け取る。ここで出力としては、「合格」または「不合格」の値が返される。この処理を文書ファイル名辞書1に登録されているすべてのファイルに対して行い、統合文書内容検査部7の出力が「合格」だったファイルに対してのみ、文書ファイル名を文書内容出力制御部10に出力する。

【0041】文書内容出力部9は、文書ファイル名とこのファイルの内容とキーワードフィルタ部5の出力である1つ以上の属性名インデックスとを入力する。入力した属性名インデックスのうちの1つを取り出し、入力した文書ファイルの内容中に、この属性名インデックス中の属性名を含むタグが存在するかどうか調べ、存在する場合は、その属性名タグに対応する属性値タグの値と入力した文書ファイル名との対を表示装置102を通じて利用者に表示し、存在しなかった場合には何も出力しないという処理を、入力したすべての属性名インデックスのそれぞれに対して行う。文書内容出力部9は、文書内容出力制御部10によってサブルーチン的に呼び出される役

割をしている。なお、属性値タグの値と文書ファイル名との対を出力する代わりに、属性値タグの値と文書ファイルの位置情報とを表示するようにしても良く、また、属性値タグの値と文書ファイル名とその位置情報とを表示するようにしても良い。

【0042】文書内容出力制御部10は、合格文書ファイル名選別部8の出力である文書ファイル名の集合をそのまま自身の入力とし、文書ファイル名辞書1を参照して、入力した文書ファイル名の集合中の文書ファイルの内容を必要に応じてインターネット103を通じてWWWサーバをアクセスして1つずつ取り出し、文書ファイル名およびキーワードフィルタ部5で生成された属性名インデックスとともに文書内容出力部9に渡すことを、入力中の文書ファイル名のすべてに対して行うことを繰り返すものである。つまり、入力として合格した文書ファイル名を3つ受け取った場合には、3回文書内容出力部9を呼び出すことになる。なお、合格文書ファイル名選別部8がインターネット103を通じてWWWサーバから取り込んだ文書ファイルの内容が磁気ディスク装置等に保存されている場合、文書内容出力制御部10はその内容を利用することで、インターネット103へのアクセス回数を減らすことができる。

【0043】図2および図3は文書ファイル検索装置100の処理例を示すフローチャートである。以下、本実施の形態の動作について説明する。

【0044】キーワード抽出部4は、入力装置101を通じて利用者から自然言語表現による検索入力文を受け付けると(ステップS1)、属性名格納辞書2と属性値格納辞書3とを参照して、その中の自然言語表現として登録されている表現が検索入力文にないかどうかを、検索入力文の先頭から順に調べ、あった場合には、それが属性名のときは属性名のみを含む属性名インデックスを出力し、それが属性値のときは属性値と対応する属性名との対を含む属性値インデックスを出力する(ステップS2)。

【0045】次にキーワードフィルタ部5は、キーワード抽出部4から出力されたインデックスの並びを検査し、同一の属性の属性名インデックスと属性値インデックスとが連続している箇所を検出し、その箇所の属性名インデックスを削除する(ステップS3)。

【0046】次に合格文書ファイル名選別部8は、文書ファイル名辞書1中の1つの文書ファイル名に注目し、その文書ファイル名の文書の内容を取り出して、キーワードフィルタ部5から出力された全ての属性値インデックスとともに統合文書内容検査部7に渡し、合否を判定させる(ステップS4)。

【0047】統合文書内容検査部7は、渡された文書内容を検査するために、まず渡された属性値インデックスの1つに注目し、この属性値インデックスと文書ファイルの内容とを文書内容検査部6に渡し、合否を判定させ

る(ステップS5)。

【0048】文書内容検査部6は、渡された文書ファイルの内容中に、渡された属性値インデックスに含まれる属性名を持つ属性名タグが存在し、かつ、その存在した属性名タグと対になっている属性値タグの値が、渡された属性値インデックスに含まれる属性値と一致するかを検査し、一致する場合には合格を、そのような属性名タグが存在しないか或いは存在してもその属性値が一致しない場合には不合格を、統合文書内容検査部7に通知する(ステップS6)。

【0049】統合文書内容検査部7は、文書内容検査部6から合格が通知された場合(ステップS7でYES)、合格文書ファイル名選別部8から通知された全ての属性値インデックスについて検査し終えたか否かを調べ、未だ検査し終えていないときは(ステップS8でNO)、残りの属性値インデックスの1つに注目を移し、その属性値インデックスと文書ファイルの内容とを文書内容検査部6に渡し、合否を判定させる(ステップS9)。そして、全ての属性値インデックスについて文書内容検査部7で合格の判定が出た場合(ステップS8でYES)、合格文書ファイル名選別部8に合格を通知し、合格文書ファイル名選別部8は当該文書ファイルを合格文書ファイルとし(ステップS10)、ステップS11へと進む。他方、文書内容検査部6から不合格が通知された場合(ステップS7でNO)、統合文書内容検査部7は合格文書ファイル名選別部8に不合格を通知し、合格文書ファイル名選別部8はステップS11へと進む。

【0050】合格文書ファイル名選別部8は、1つの文書ファイルについての合否判定が終わると、文書ファイル名辞書1中に未処理の文書ファイルが残っている場合(ステップS11でYES)、その内の1つの文書ファイル名に注目を移し(ステップS12)、先の文書ファイルと同様に合否の判定を下す。

【0051】文書ファイル名辞書1中の全ての文書ファイルに対する合否判定を終えると(ステップS11でYES)、合格文書ファイル名選別部8は、少なくとも1つの合格ファイルがあったか否かを判定し(ステップS13)、1つもなければ、例えば入力された検索条件に合致する文書ファイルは1つもなかった旨を利用者に表示する等の処理を行って、処理を終了する。他方、1つでも合格ファイルが存在した場合、その全ての合格ファイルの文書ファイル名とキーワードフィルタ部5から出力された全ての属性名インデックスとを文書内容出力制御部10に通知して、文書内容出力制御を開始させる(ステップS14)。

【0052】文書内容出力制御部10は、通知された1つの合格ファイル名に注目してその文書内容を取り出し、通知された全ての属性名インデックスとともに文書内容出力部9に渡し、当該文書の処理を開始させる(ステップS15)。

【0053】文書内容出力部9は、通知された1つの属性名インデックスに注目し（ステップS16）、その属性名インデックスの属性名をもつ属性名タグが文書内にあるかを調べ（ステップS17）、あれば（ステップS18でYES）、その属性名タグに対応する属性値タグの値と当該文書ファイル名とを表示装置102に表示する（ステップS19）。なければ（ステップS18でNO）、ステップS19をスキップする。次に文書内容出力部9は、通知された属性名インデックスに未処理の属性名インデックスが残っているか否かを調べ（ステップS20）、残っていれば、その1つに注目を移し（ステップS21）、ステップS17に戻って上述した処理を繰り返す。

【0054】文書内容出力部9が通知された全ての属性名インデックスについての処理を終えると（ステップS20でNO）、文書内容出力制御部10は、合格文書ファイル名選別部8から通知された文書ファイルに未処理のものが残っているか否かを調べ（ステップS22）、残っている場合にはその1つに注目を移し、その文書ファイル名の文書内容を取り出して、合格文書ファイル名選別部8から通知された全ての属性名インデックスとともに文書内容出力部9に渡し、処理させる（ステップS23）。全ての合格ファイルについての処理が終わると（ステップS22でYES）、処理終了となる。

【0055】

【実施例】文書ファイル名辞書1に、図1に例示するように「ファイル1」、「ファイル2」、「ファイル3」の3つの文書ファイル名とそのURLとが登録されているとする。また、ファイル1の内容が図5(a)に示すものであり、ファイル2の内容が図5(b)に示すものであるとする。これらのファイル1、2はXMLで記述されており、文章テキストの他に属性と属性値が含まれている。つまり、ファイル1には、＜ホテル＞Xホテル／ホテル＞、＜場所＞サンフランシスコ／場所＞、＜値段＞\$150／値段＞、＜ベット＞可能／ベット＞といった、文書中に書かれた意味を表現する属性名のついたタグとその属性の値との対が含まれている。同様に、ファイル2にも、＜ホテル＞Zホテル／ホテル＞、＜場所＞シアトル／場所＞、＜値段＞\$180／値段＞、＜ベット＞不可／ベット＞といったタグが含まれている。

【0056】また、属性名格納辞書2には図1に例示するような属性名とその自然言語表現との対が事前に格納されており、属性値格納辞書3には図1に例示するような属性値と自然言語表現と属性名との3つ組が事前に格納されているものとする。なお、属性値格納辞書3に全ての価格をその実際値で登録すると、登録数が増えてしまうので、変数を使用して登録するようにしても良い。つまり、XXXを任意の数値とする場合、以下の表3に示すように登録しておき、キーワード抽出部4は任意の数値の後ろに「ドル」があれば、自然言語表現XXXドルが

存在すると判断し、存在した実際値の頭に\$を付けたものを属性値とする。

【表3】

属性値	自然言語表現	属性名
\$XXX	XXXドル	価格

【0057】このような前提で、利用者が以下のような自然言語による検索入力文を入力した場合を例に、本実施例の動作を説明する。

検索入力文：「値段が150ドルでペットを連れて行くのが可能で、しかも、サンフランシスコ近辺にあるホテルの情報を見つけたい」

【0058】キーワード抽出部4は利用者からの検索入力文を受け付けると、属性名格納辞書2および属性値格納辞書3を参照して、検索入力文を以下のようにキーワード列に変換する。

【0059】まず、検索入力文の先頭の自然言語表現「値段」が属性名格納辞書2に存在するので、それと対になって登録されている属性名「価格」を属性名インデックスとして出力する。次に、自然言語表現「150ドル」が属性値格納辞書3に存在するので、それと3つ組で登録されている属性値「\$150」と属性名「価格」との対を属性値インデックスとして出力する。次に、自然言語表現「ペット」が属性名格納辞書2に存在するので、それと対になって登録されている属性名「ペット」を属性名インデックスとして出力する。次に、自然言語表現「可能」が属性値格納辞書3に存在するので、それと3つ組で登録されている属性値「可能」と属性名「ペット」との対を属性値インデックスとして出力する。次に、自然言語表現「サンフランシスコ」が属性値格納辞書3に存在するので、それと3つ組で登録されている属性値「サンフランシスコ」と属性名「場所」との対を属性値インデックスとして出力する。次に、「ホテル」が属性名格納辞書2に存在するので、それと対になって登録されている属性名「ホテル」を属性名インデックスとして出力する。検索入力文中には、属性名格納辞書2および属性値格納辞書3に登録された自然言語表現とマッチする他の自然言語表現はない。従って、以下のようなキーワード列が上から順に出力される。

【0060】属性名インデックス（属性名「価格」）
属性値インデックス（属性値「\$150」、属性名「価格」）

属性名インデックス（属性名「ペット」）

属性値インデックス（属性値「可能」、属性名「ペット」）

属性値インデックス（属性値「サンフランシスコ」、属性名「場所」）

属性名インデックス（属性名「ホテル」）

【0061】次にキーワードフィルタ部5は、属性名インデックスと属性値インデックスとの並び順を参照して、冗長な部分の統合を行う。上のキーワード列の場合、属性名インデックス（属性名「価格」）と属性値インデックス（属性値「\$150」、属性名「価格」）とは同じ属性名「価格」で隣どうしに並んでいるので、属性名インデックス（属性名「価格」）を削除する。また、属性名インデックス（属性名「ベット」）と属性値インデックス（属性値「可能」、属性名「ベット」）とは同じ属性名「ベット」で隣どうしに並んでいるので、属性名インデックス（属性名「ベット」）を削除する。他に削除すべき属性名インデックスは存在しないので、上記のキーワード列は最終的に以下のように圧縮される。

- 【0062】(a) 属性値インデックス（属性値「\$150」、属性名「価格」）
 (b) 属性値インデックス（属性値「可能」、属性名「ベット」）
 (c) 属性値インデックス（属性値「サンフランシスコ」、属性名「場所」）
 (d) 属性名インデックス（属性名「ホテル」）

【0063】次に、合格文書ファイル名選別部8は、文書ファイル名辞書1中のファイル1の文書内容をそのURLを頼りにインターネット103を通じて該当するサーバから取得し、その文書内容と上記の属性値インデックス(a)～(c)とを統合文書内容検査部7に渡す。

【0064】統合文書内容検査部7は、ファイル1の文書内容と、1つの属性値インデックス(a)とを文書内容検査部6に渡す。

【0065】文書内容検査部6は、ファイル1の文書内容中に、属性値インデックス(a)中の属性名「価格」のタグが存在するか否かを調べる。図5(a)のファイル1の場合、該当するタグ<値段>\$150</値段>があるので、その属性値「\$150」が受け取った属性値インデックス(a)中の属性値「\$150」と一致するか否かを調べる。今の例では、一致するので、合格を統合文書内容検査部7に返却する。

【0066】統合文書内容検査部7は、ファイル1の文書内容と、次の属性値インデックス(b)とを文書内容検査部6に渡す。

【0067】文書内容検査部6は、ファイル1の文書内容中に、属性値インデックス(b)中の属性名「ベット」のタグが存在するか否かを調べる。図5(a)のファイル1の場合、該当するタグ<ベット>可能</ベット>があるので、その属性値「可能」が受け取った属性値インデックス(b)中の属性値「可能」と一致するか否かを調べる。今の例では、一致するので、合格を統合文書内容検査部7に返却する。

【0068】統合文書内容検査部7は、ファイル1の文書内容と、次の属性値インデックス(c)とを文書内容検査部6に渡す。

【0069】文書内容検査部6は、ファイル1の文書内容中に、属性値インデックス(c)中の属性名「場所」のタグが存在するか否かを調べる。図5(a)のファイル1の場合、該当するタグ<場所>サンフランシスコ</場所>があるので、その属性値「サンフランシスコ」が受け取った属性値インデックス(c)中の属性値「サンフランシスコ」と一致するか否かを調べる。今の例では、一致するので、合格を統合文書内容検査部7に返却する。

10 【0070】統合文書内容検査部7は、ファイル1に関し全ての属性値インデックスで合格の結果が得られたので、合格文書ファイル名選別部8に合格を通知し、合格文書ファイル名選別部8はファイル1を合格ファイルとする。

【0071】次に合格文書ファイル名選別部8は、文書ファイル名辞書1に格納されたファイル2の文書内容をそのURLを頼りにインターネット103を通じて該当するサーバから取り込み、先のファイル1と同様に統合文書内容検査部7を使って合否を判定する。この場合、ベット属性、場所属性が満足しないので、ファイル2は不合格となる。同様に、残りのファイル3についても合否の判定が行われる。ここでは、ファイル3も不合格と判定され、合格ファイルはファイル1のみであったとする。

【0072】次に合格文書ファイル名選別部8は、合格ファイル名としてファイル名1を、属性名インデックス(d)とともに文書内容出力制御部10に渡す。

30 【0073】文書内容出力制御部10は、文書ファイル名辞書1からファイル名1のURLを取得し、それを頼りにインターネット103上のサーバをアクセスしてファイル名1の文書内容を取得し、属性名インデックス(d)とともに文書内容出力部9に渡す。

【0074】文書内容出力部9は、ファイル1の文書内容中に、属性名インデックス(d)の属性名「ホテル」を持つ属性タグが存在するか否かを調べる。図5(a)のファイル1の場合、該当するタグ<ホテル>Xホテル</ホテル>が存在するので、その属性値「Xホテル」を取り出し、ファイル名1と共に表示装置102に表示する。

40 【0075】図4は本発明の文書ファイル検索装置のハードウェア構成例を示すブロック図である。この例の文書ファイル検索装置は、プロセッサ(CPU)200と、磁気ディスク装置等の補助記憶装置201と、そのインタフェース202と、RAM等のメモリ203と、インターネット103との間のインタフェース204と、CD-ROM、半導体メモリ等の機械読み取り可能な記録媒体205と、そのインタフェース206と、入力装置101と、そのインタフェース207と、表示装置102と、そのインタフェース208と、CPU200、メモリ203およびインタフェース202、20

25

4、206～208間を接続するバス209とから構成されている。

【0076】記録媒体205には、文書ファイル検索用プログラムが記録されており、このプログラムがインタフェース206を介してインストールされることにより、メモリ203または補助記憶装置201上に図1の文書ファイル名辞書1、属性名格納辞書2および属性値格納辞書3がロードされる。また、同プログラムはCPU200の動作を制御することにより、CPU200を図1のキーワード抽出部4、キーワードフィルタ部5、文書内容検査部6、統合文書内容検査部7、合格文書ファイル名選別部8、文書内容出力部9、文書内容出力制御部10として機能させる。

【0077】以上の実施の形態は本発明をWWW上のホームページの検索に適用したが、特開平5-67136号公報に記載する技術と同様にデータベースに対する検索に対しても適用可能である。

【0078】

【発明の効果】以上説明したように本発明によれば以下のような効果が得られる。

【0079】自然言語による検索問い合わせに対する回答の冗長性を極力無くすることができる。その理由は、キーワード抽出部で抽出されたインデックス列をその先頭から順に探査し、同一の属性の属性名インデックスと属性値インデックスとが隣どうしに存在する場合に属性名インデックスを削除するキーワードフィルタ部を備えているからである。

【0080】WWWのホームページに対しても自然言語による検索問い合わせが可能になる。その理由は、XMLの

26

ように文書ファイル中にその意味を表現する属性名と属性値のタグを内蔵させており、利用者が入力した自然言語による検索条件を解釈して適合する属性名および属性値をもつ文書ファイル中から、利用者の望む属性値を取り出すことができるからである。これによって、文法的に正しい自然言語表現、非文法的な表現、自然言語文の断片、キーワード列等、種々の形の入力を受け付けて統一的に解釈を行なう自然言語インタフェースによるWWW文書検索システムを実現することが出来る。

【図面の簡単な説明】

【図1】本発明の実施の形態の文書ファイル検索装置のブロック図である。

【図2】本発明の実施の形態の文書ファイル検索装置処理例を示すフローチャートである。

【図3】本発明の実施の形態の文書ファイル検索装置処理例を示すフローチャートである。

【図4】本発明の文書ファイル検索装置のハードウェア構成例を示すブロック図である。

【図5】XMLを使った文書ファイルの記述例を示す図である。

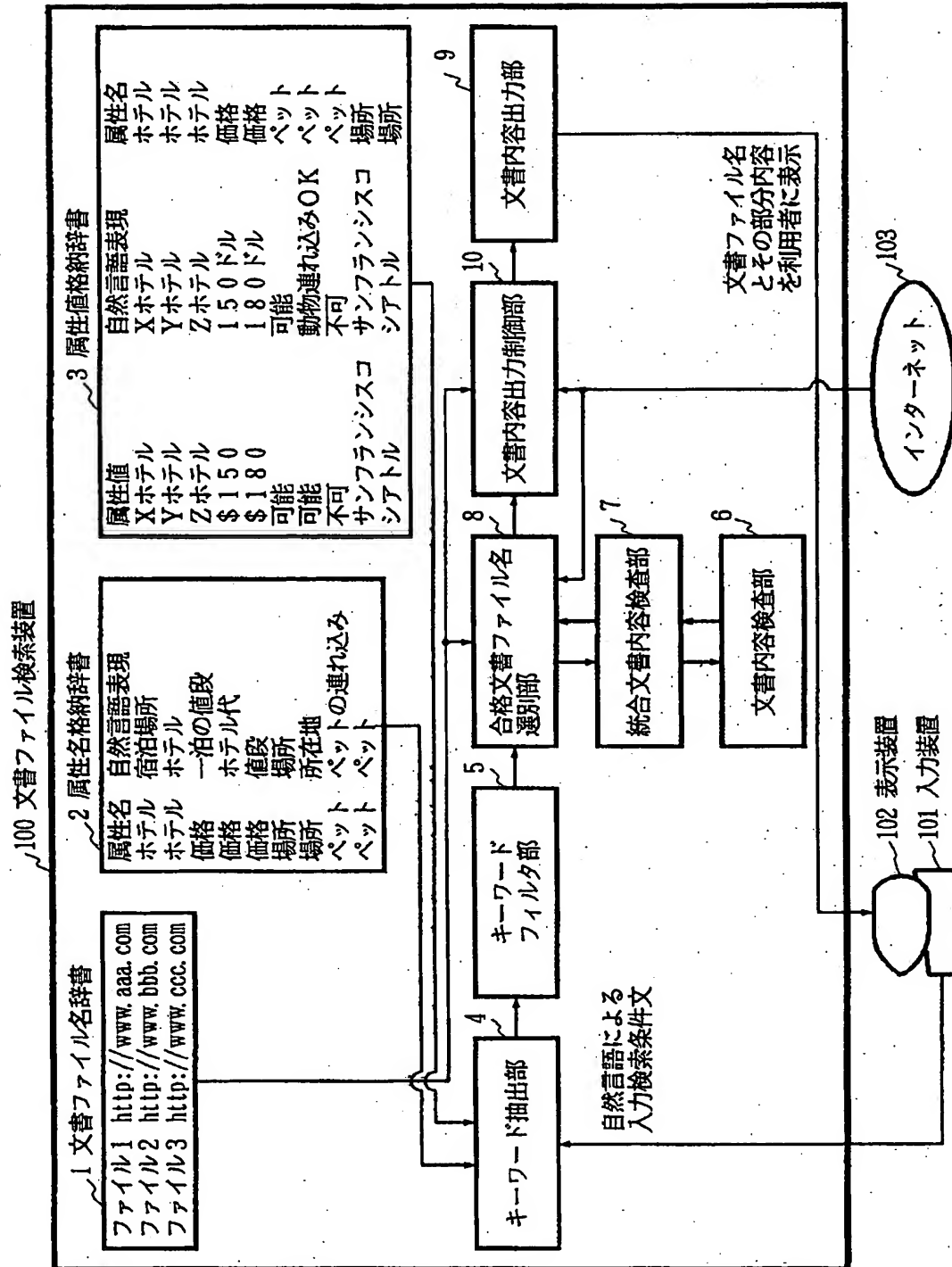
【図6】HTMLの記述例を示す図である。

【図7】図6のHTMLの記述例をブラウザで表示した例を示す図である。

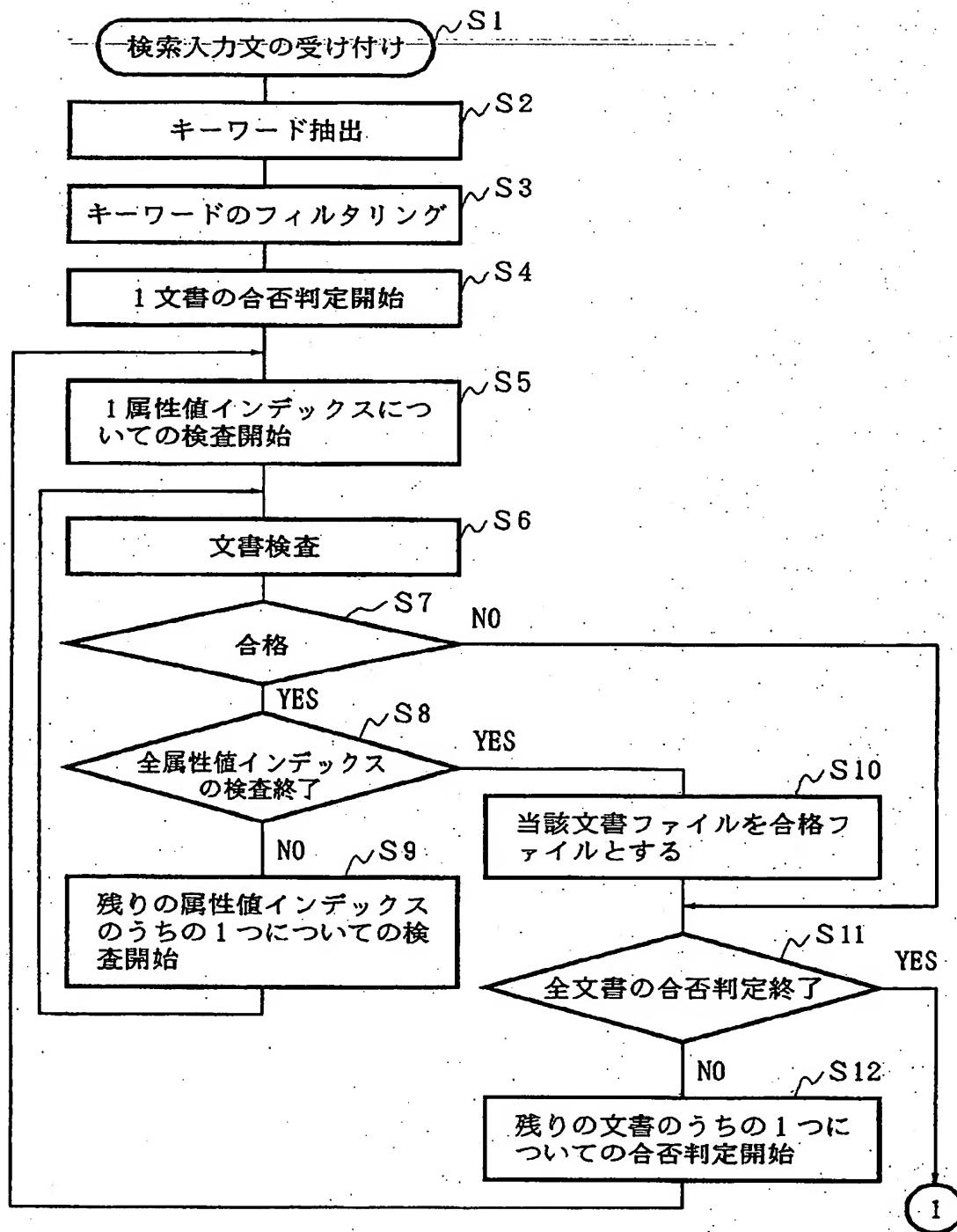
【符号の説明】

1は文書ファイル名辞書、2は属性名格納辞書、3は属性値格納辞書、4はキーワード抽出部、5はキーワードフィルタ部、6は文書内容検査部、7は統合文書内容検査部、8は合格文書ファイル名選別部、9は文書内容出力部、10は文書内容出力制御部、である。

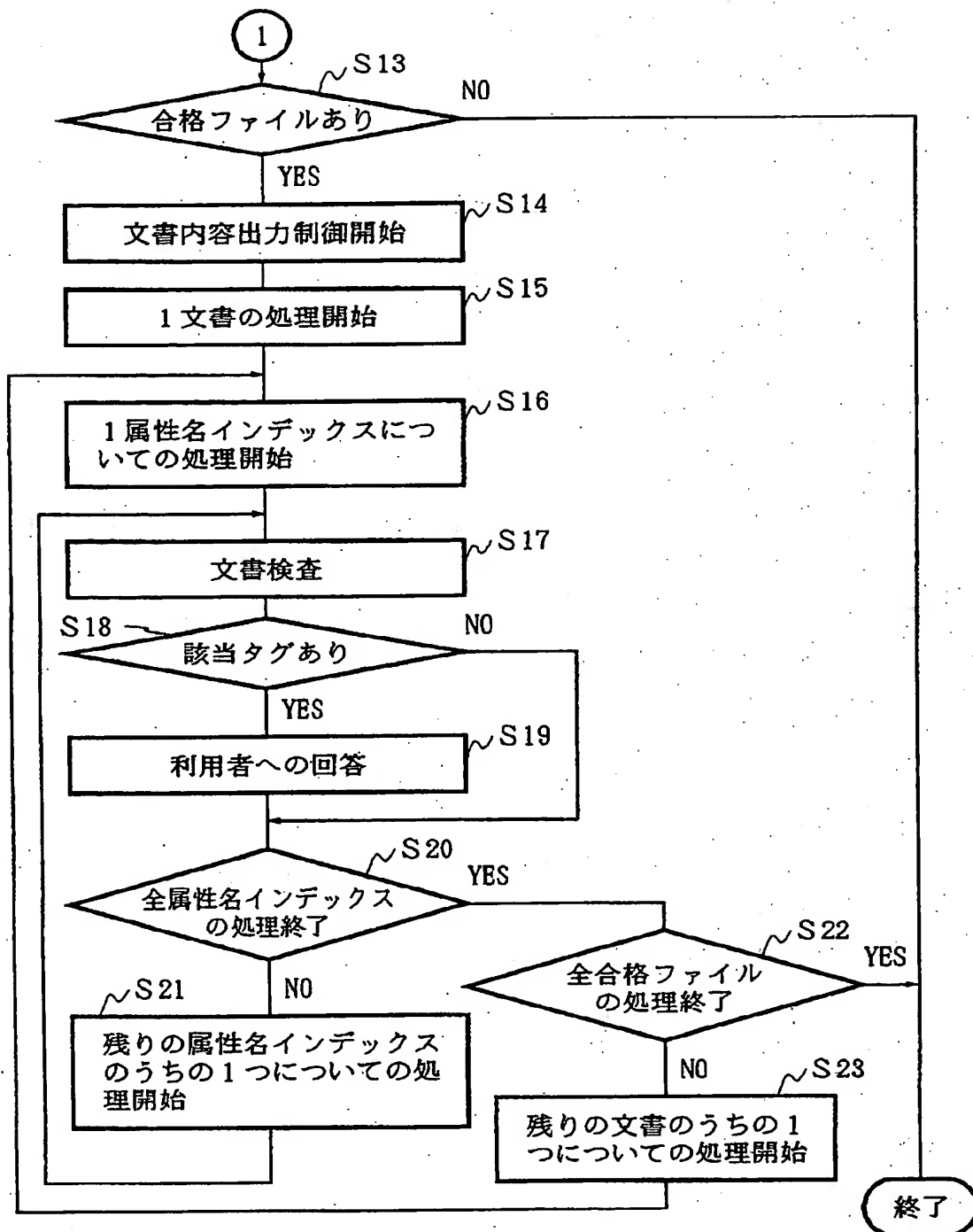
【図 1】



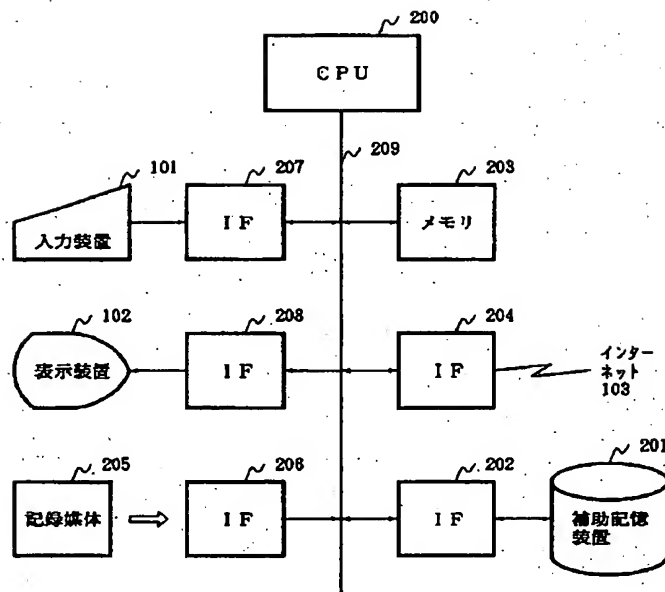
【図2】



【図3】



【図4】



【図 5】

(a)

```
<!--Xホテルのプロフィールのファイル--> <====これはコメント行
<?xml version="1.0", encoding="shiftjis"?>
<!DOCTYPE 宿泊情報プロフィールSYSTEM"http://report/hotel.xml">
<プロフィール記録>
<ホテル>Xホテル</ホテル>
<場所>サンフランシスコ</場所>
<値段>$150</値段>
<ペット>可能</ペット>

<プロフィール>
ホテルXは、ビジネス街に位置し、ビジネス、観光いずれにも、...
</プロフィール>
...
</プロフィール記録>
```

(b)

```
<!--Zホテルのプロフィールのファイル--> <====これはコメント行
<?xml version="1.0", encoding="shiftjis"?>
<!DOCTYPE 宿泊情報プロフィールSYSTEM"http://report/hotel.xml">
<プロフィール記録>
<ホテル>Zホテル</ホテル>
<場所>シアトル</場所>
<値段>$180</値段>
<ペット>不可</ペット>

<プロフィール>
ウォーターフロントの中心部にあるZホテルは、シアトルで最もホットな場所として、...
</プロフィール>
...
</プロフィール記録>
```

【図 6】

```

<html><title>1201</title>
<body BGCOLOR="#ffffff" TEXT="#000000">
<TABLE BORDER="1" CELSPACING="1" WIDTH="550">
<TR><TH>担当課題名</TH><TD COLSPAN="2">XMLの自然言語インタフェース方式の開発</TD></TR>
<TR><TH>担当名</TH><TD>N E C   C & Cメディア研究所</TD></TR>
<TR><TH>優先度</TH><TD>高い</TD></TR>
</TABLE>
<P>
<B>1. 研究目的</B><P>
本研究は、先端的であり、....
<P>
<B>2. 具体的既往成果</B><P>
本研究部門では、すでに、自然言語インタフェースの研究を、....
<P>
<B>3. 全体計画</B><P>
1年目には、....、2年目には、....
</body></HTML>

```


【図7】

担当課題名 XMLの自然言語インタフェース方式の開発

担当者 NEC C&Cメディア研究所

優先度 高い

1. 研究目的

本研究は、先端的であり、……

2. 具体的既往成果

本研究部門では、すでに、自然言語インタフェースの研究を、……

3. 全体計画

1年目には、……、2年目には、……